

PROJET 2

ANALYSEZ DES DONNÉES DE SYSTÈMES ÉDUCATIFS



PLAN

01 RAPPEL DE LA PROBLÉMATIQUE ET PRÉSENTATION DU JEU DE DONNÉES

Contexte, problématique, objectifs et description du dataset

02 ANALYSE PRÉ-EXPLORATOIRE

Environnement de développement, Analyse de la qualité, Sélection des données pertinentes, Scoring et Visualisation des résultats

03 CONCLUSION ET PERSPECTIVES

Suggestions de réponses aux questions métier et perspectives d'amélioration



01

PROBLÉMATIQUE ET JEU DE DONNÉES

Contexte, problématique, objectifs
et description du dataset

CONTEXTE

- academy : start-up de la EdTech
- Métier: propose des **contenus de formation en ligne**
- Cible: public de niveau **lycée** et **université**
- Projet d'**expansion** à l'**international**
- Identification du **potentiel de nouveaux marchés** (en explorant un jeu de données de la Banque Mondiale sur le thème de l'éducation)





QUESTIONS

- Quels sont les pays avec un fort potentiel de clients pour les services de l'entreprise ?
- Pour chacun de ces pays, quelle sera l'évolution de ce potentiel de clients ?
- Dans quels pays l'entreprise doit-elle opérer en priorité ?



VALIDATION DE LA QUALITÉ DU JEU DE DONNÉES

Analyser la qualité du jeu de données (comporte-t-il beaucoup de données manquantes, dupliquées ?)



DESCRIPTION DU CONTENU DU JEU DE DONNÉES

Décrire les informations contenues dans le jeu de données (nombre de colonnes ? nombre de lignes ?)

OBJECTIFS DE LA MISSION



SELECTION DES DONNÉES PERTINENTES

Sélectionner les colonnes contenant des informations qui peuvent être utiles pour répondre à la problématique de l'entreprise



ANALYSE PRÉ- EXPLORATOIRE

Déterminer des ordres de grandeurs des indicateurs statistiques classiques pour les différentes zones géographiques et pays du monde

SOURCE DU JEU DE DONNÉES



THE WORLD BANK
IBRD • IDA

Collectées par l'organisme "EdStats All Indicator Query" de la
Banque mondiale sur des sujets liés à l'éducation



INDICATEURS

3665 indicateurs
internationaux

SUJETS

Indicateurs répartis
en 37 sujets

RÉGIONS

7 régions du globe

PAYS

242 pays et zones

1970-2100

DESCRIPTION

FICHER	#LIGNES	#COLONNES	DESCRIPTION
EdStatsData.csv	886930	70	Valeurs des indicateurs pour tous les pays (1970-2100)
EdStatsSeries.csv	3665	21	Infos sur les indicateurs regroupés par sujets (topics)
EdStatsCountry.csv	241	32	Infos globales sur l'économie de chaque pays/zone du monde
EdStatsCounty-Series.csv	613	4	Description des sources des indicateurs
EdStatsFootnote.csv	643638	5	Infos supplémentaires indicateurs / pays / année

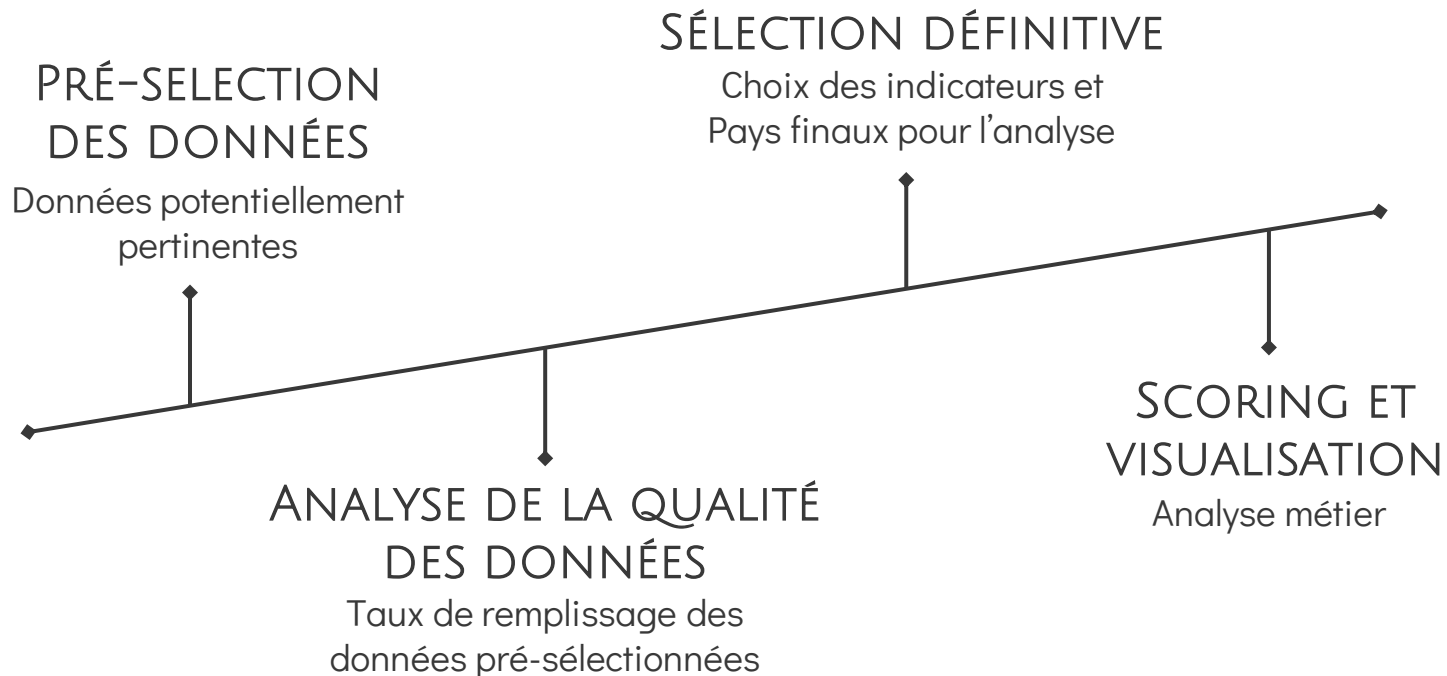


02

ANALYSE PRÉ- EXPLORATOIRE

Environnement de développement,
Analyse de la qualité, Sélection des
données pertinentes, Scoring et
Visualisation des résultats

STRATÉGIE



ENVIRONNEMENT DE DÉVELOPPEMENT

ANACONDA

Installation d'Anaconda:
plateforme de distribution
python la plus populaire

ENVIRONNEMENT VIRTUEL

Mise en place d'un
environnement virtuel dédié
au projet

INSTALLATION DES PAQUETS

Installation des paquets
nécessaires (numpy,
pandas, matplotlib, seaborn)
avec la commande **pip install**

PRÉ-SÉLECTION DES DONNÉES

PRÉ-SÉLECTION DES DONNÉES

RAPPELS

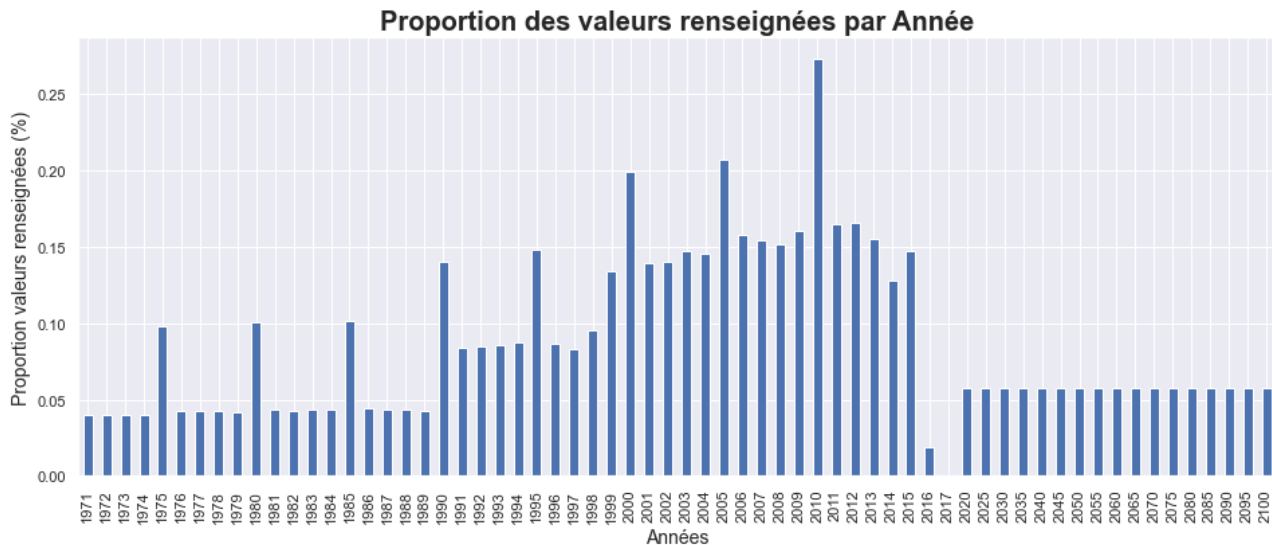
- Métier: formation en ligne
- Cible: lycée, université
- Expansion à l'international

FICHIERS IMPORTANTS

- EdStatsData.csv
- EdStatsSeries.csv
- EdStatsCountry.csv

PÉRIODE PERTINENTE

- Les données sont plus renseignées pour 2000-2015
- Ignorer les données 1970-1999



TOPIC PERTINENTS

TOPIC	DESCRIPTION	NB INDIC.
Attainment	Niveau d'instruction	733
Economic Policy & Debt: National accounts: US at current prices: Aggregate indicators	Pouvoir d'achat	3
Infrastructure: Communications	Utilisation d'internet	2
Health: Population: Dynamics	Croissance de la population	1
Health: Population: Structure	Structure de la population	13
Population	Structure population / niveau scolaire	213
Secondary	Niveau lycée	256
Tertiary	Niveau université	158
TOTAL		1379

RÉPARTITION DU DATASET

▲ 7 RÉGIONS

▲ 242 PAYS

▲ 1379 INDICATEURS

▲ 8 SUJETS (TOPIC)

▲ ANNEES 2000-2100

HISTORIQUE

2000-2015

PROJECTION

2020-2100

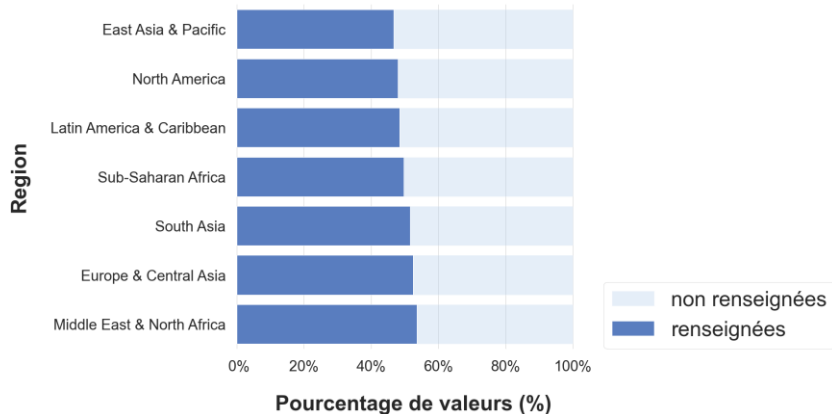
Données regroupées en 2 bases

- Les valeurs des indicateurs
- Par pays et région
- Sur les périodes

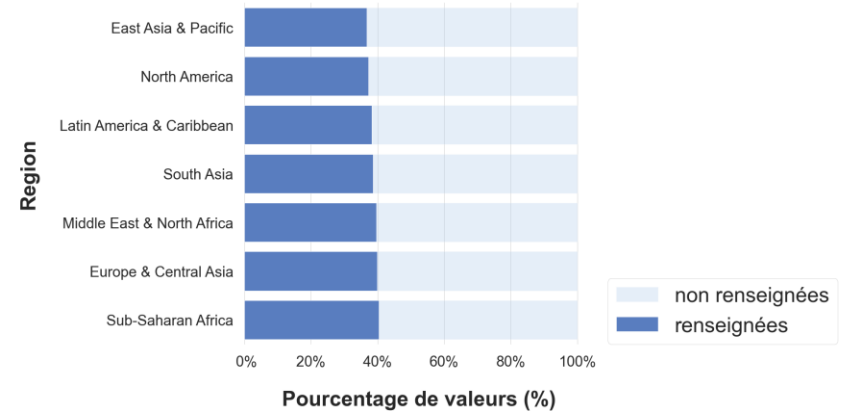
ANALYSE DE LA QUALITÉ DES DONNÉES PRÉSÉLECTIONNÉES

TAUX DE REMPLISSAGE PAR RÉGION

DONNÉES HISTORIQUES: TAUX DE REMPLISSAGE DES DONNEES PAR REGION



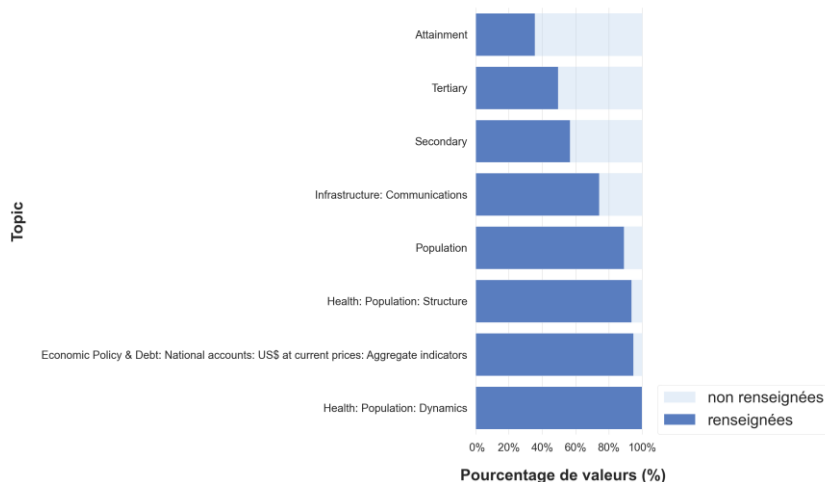
DONNÉES PROJECTION: TAUX DE REMPLISSAGE DES DONNEES PAR REGION



Données Historiques: taux de remplissage/région < 55%
Données de projection: taux de remplissage/région <= 40%

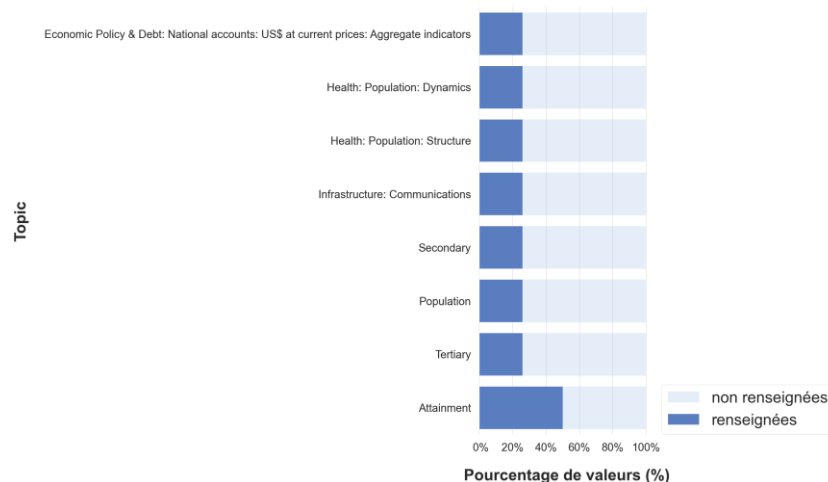
TAUX DE REMPLISSAGE PAR TOPIC

DONNÉES HISTORIQUES: TAUX DE REMPLISSAGE DES DONNÉES PAR TOPIC



- Taux de remplissage > 50% pour la plupart des topics pré-sélectionnés
- Secondary et tertiary remplis (pourtant c'est la cible)

DONNÉES PROJECTION: TAUX DE REMPLISSAGE DES DONNÉES PAR TOPIC

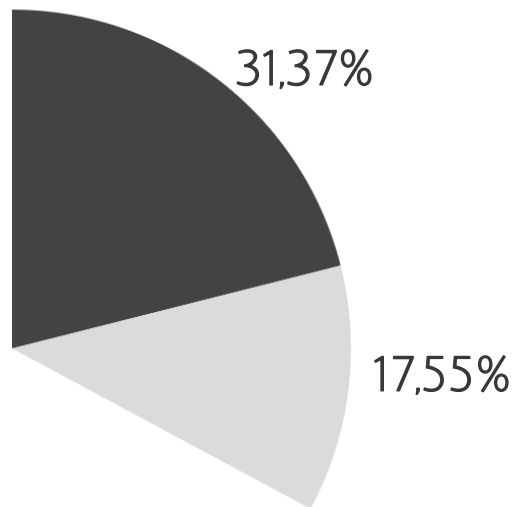


- Le topic Attainment est apparemment plus liés aux projections
- Les autres topics sont très peu remplis

TAUX DE REMPLISSAGE GLOBAL

▲ HISTORIQUE
2000-2015

▲ PROJECTION
2020-2100



REMARQUE

Jeu de données très peu rempli dans l'ensemble

Aucune ligne dupliquée dans le jeu de donnée

SÉLECTION DÉFINITIVE DES INDICATEURS ET PAYS

INDICATEURS DE DESCRIPTION

CODE INDICATEUR	TOPIC	SIGNIFICATION
NY.GDP.PCAP.CD	Economic Policy & Debt:.....	GDP per capita (current US)
IT.NET.USER.P2	Infrastructure: Communications	Internet users (per 100 people)
SP.POP.GROW	Health: Population: Dynamics	Population growth (annual %)
SP.POP.TOTL	Health: Population: Structure	Population, total
SP.POP.1564.TO.ZS	Health: Population: Structure	Population, ages 15-64 (% of total)
SP.SEC.TOTL.IN	Population	Population of the official age for secondary education, both sexes (number)
SP.TER.TOTL.IN	Population	Population of the official age for tertiary education, both sexes (number)
SE.SEC.ENRL	Secondary	Enrolment in secondary education, both sexes (number)
SE.SEC.ENRR	Secondary	Gross enrolment ratio, secondary, both sexes (%)
SE.TER.ENRL	Tertiary	Enrolment in tertiary education, both sexes (number)
SE.TER.ENRR	Tertiary	Gross enrolment ratio, tertiary, both sexes (%)

INDICATEURS DE PROJECTION

CODE INDICATEUR	TOPIC	SIGNIFICATION
PRJ.ATT.15UP.2.MF	Attainment	Wittgenstein Projection: Percentage of the population age 15+ by highest level of educational attainment. Lower Secondary. Total
PRJ.ATT.15UP.3.MF	Attainment	Wittgenstein Projection: Percentage of the population age 15+ by highest level of educational attainment. Upper Secondary. Total
PRJ.ATT.15UP.4.MF	Attainment	Wittgenstein Projection: Percentage of the population age 15+ by highest level of educational attainment. Post Secondary. Total

RÉPARTITION DES INDICATEURS

■ INTERNET

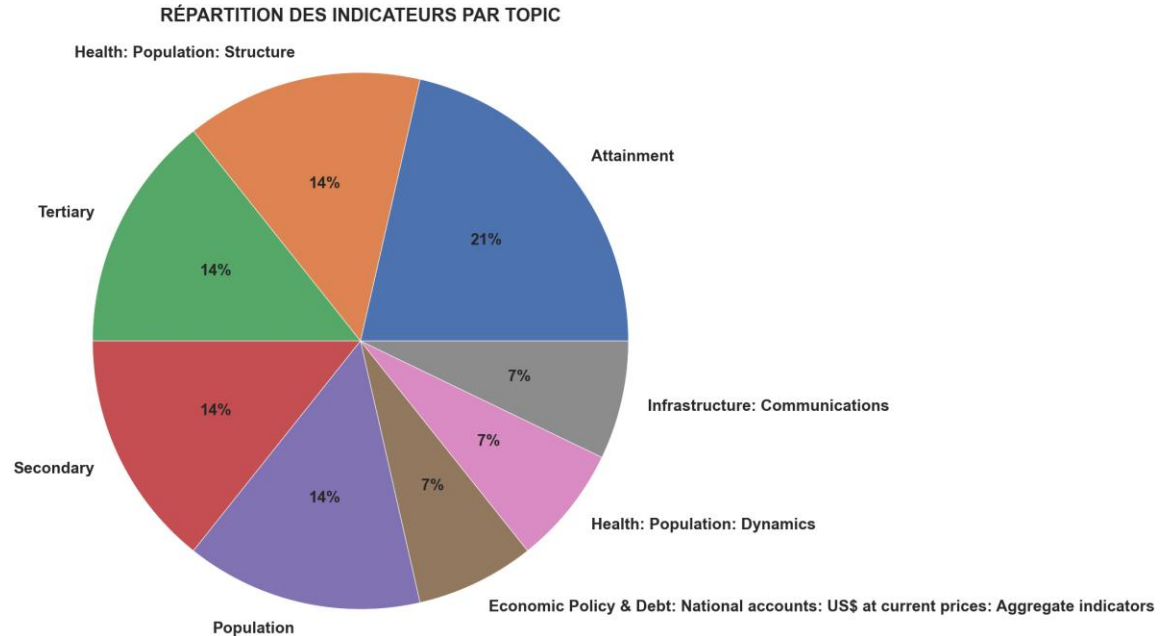
■ NIVEAU LYCÉE

■ NIVEAU UNIVERSITÉ

■ DÉMOGRAPHIE

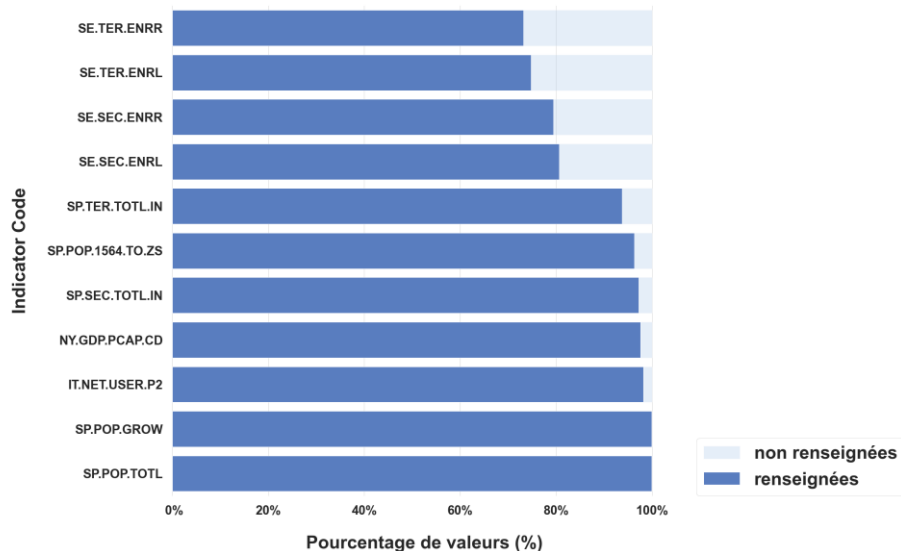
■ ÉCONOMIE

■ PROJECTION

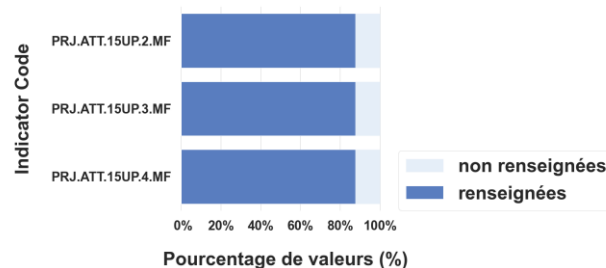


TAUX DE REMPLISSAGE DES INDICATEURS DÉFINITIFS

DONNÉES HISTORIQUES: TAUX DE REMPLISSAGE DES DONNEES PAR INDICATOR CODE



DONNÉES PROJECTION: TAUX DE REMPLISSAGE DES DONNEES PAR INDICATOR CODE



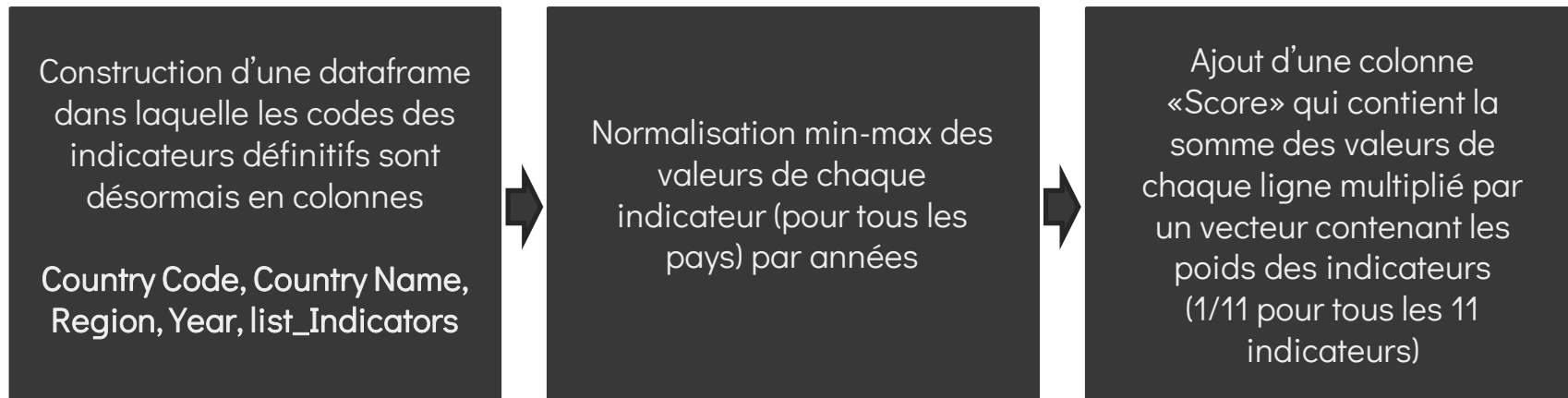
SELECTION DES PAYS

Suppression de tous les pays ayant
moins de **30%** de taux de
remplissage

SCORING ET VISUALISATION

SCORING

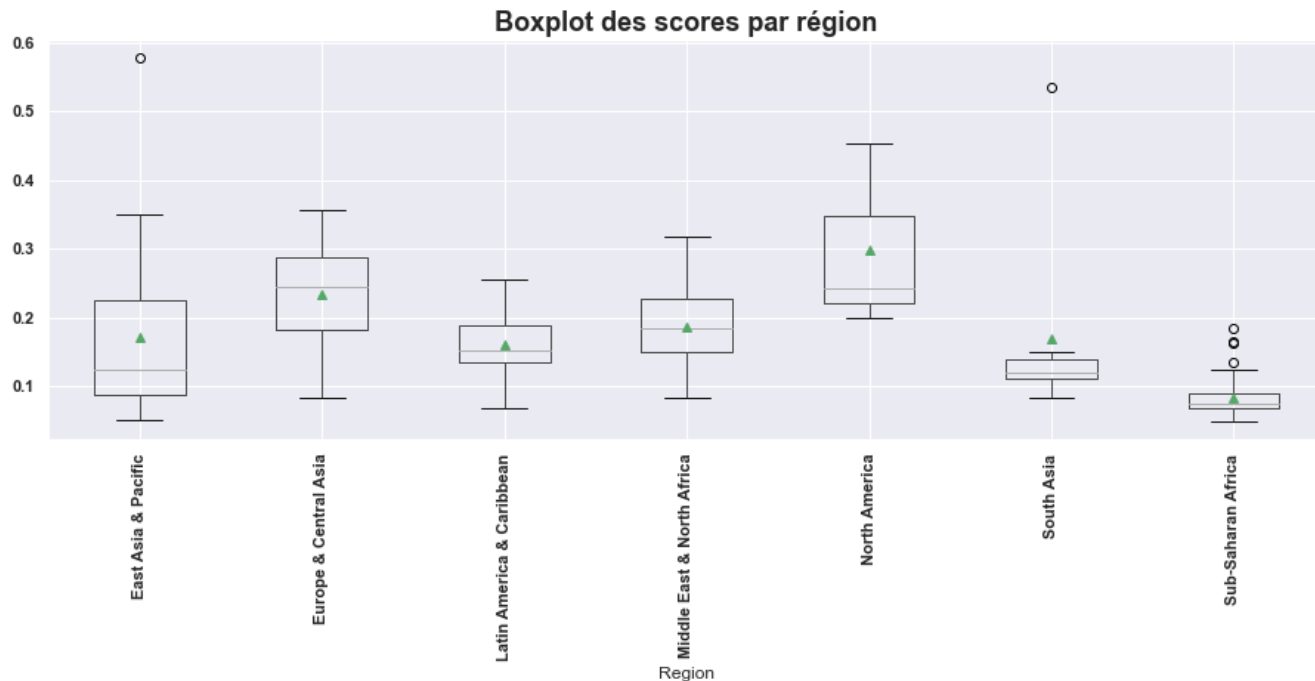
Objectif: attribuer un score à chaque Pays



Interprétation du score

Plus le score d'un pays est élevé plus son potentiel de clients pour l'entreprise est grand

LES SCORES PAR RÉGION



POTENTIEL

Les régions North America et Europe & Central Asia semblent se démarquer des autres régions

OUTLIERS

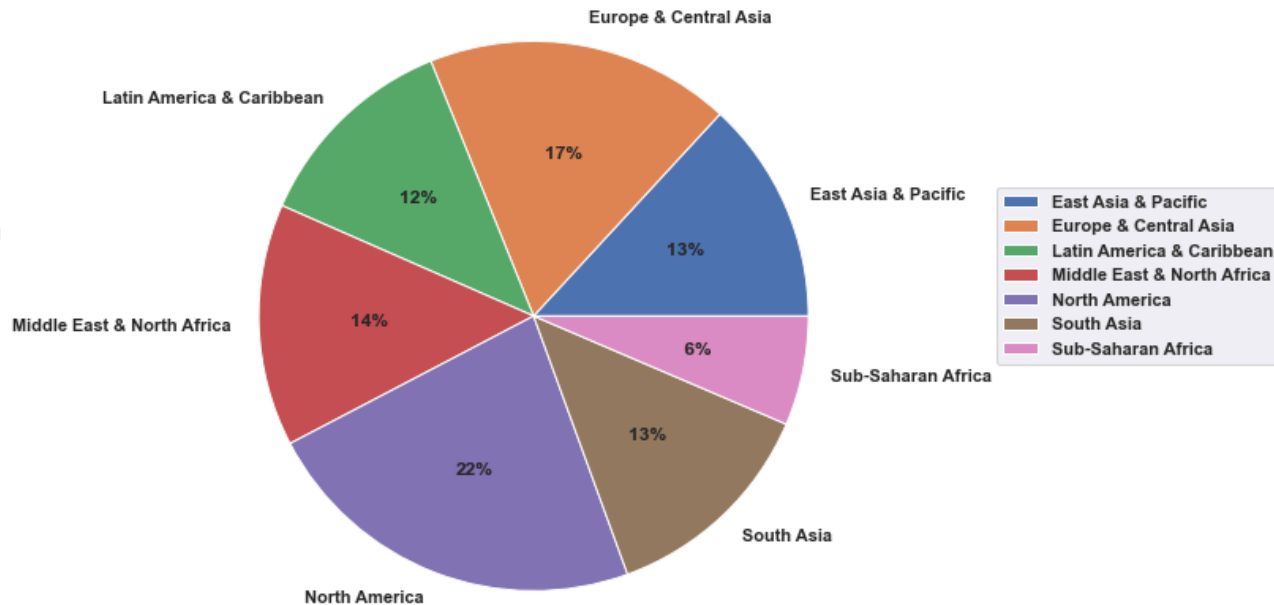
On observe des outliers (pays avec un score plus élevé par rapport aux autres pays de la région) pour les régions East Asia & Pacific, South Asia et Sub-Saharan Africa

POTENTIEL PAR RÉGION (SCORE)

REPARTITION DU POTENTIEL PAR REGION

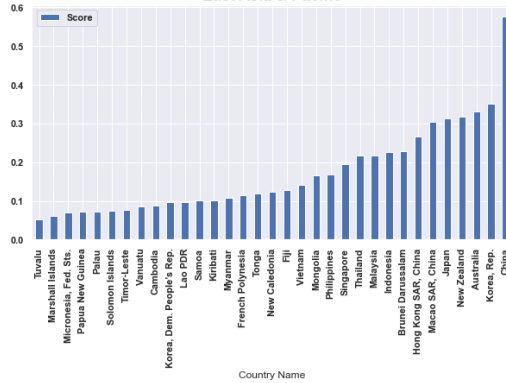
POTENTIEL

L'hypothèse émise avec le boxplot du score par région sur la démarcation des régions North America et Europe & central Asia se confirme

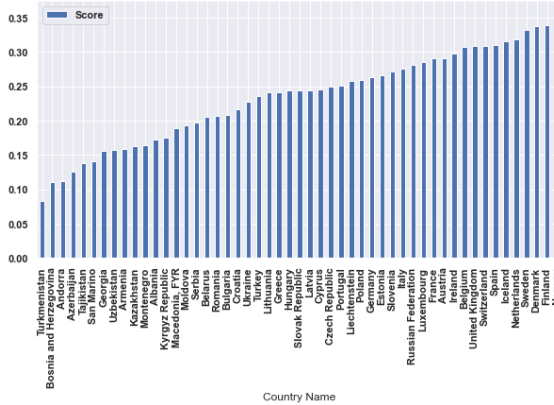


SCORES PAR PAYS

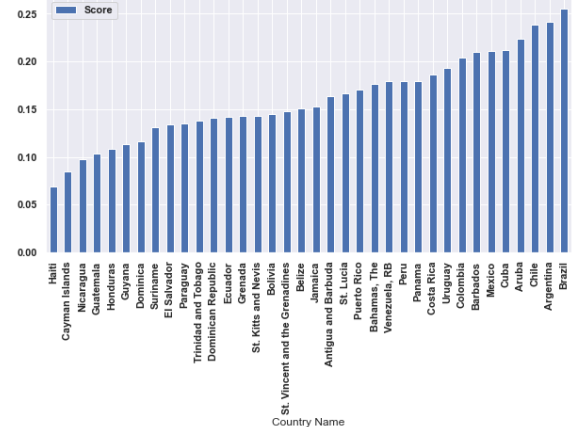
East Asia & Pacific



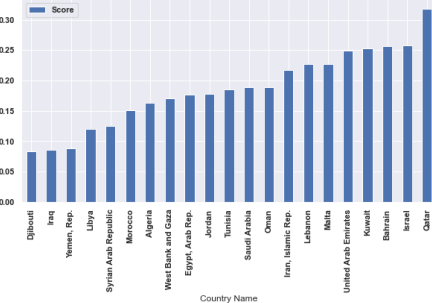
Europe & Central Asia



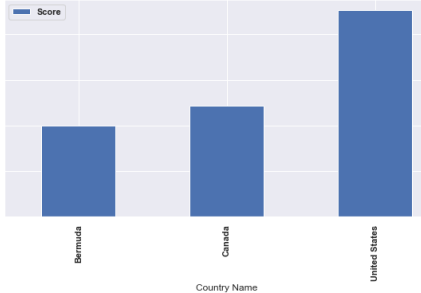
Latin America & Caribbean



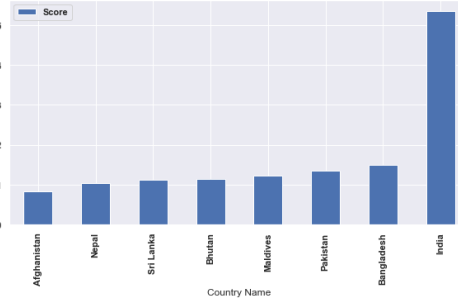
Middle East & North Africa



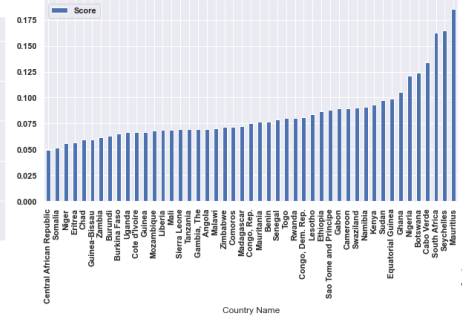
North America



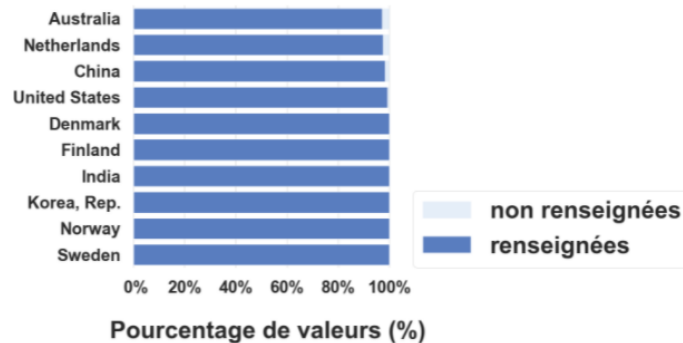
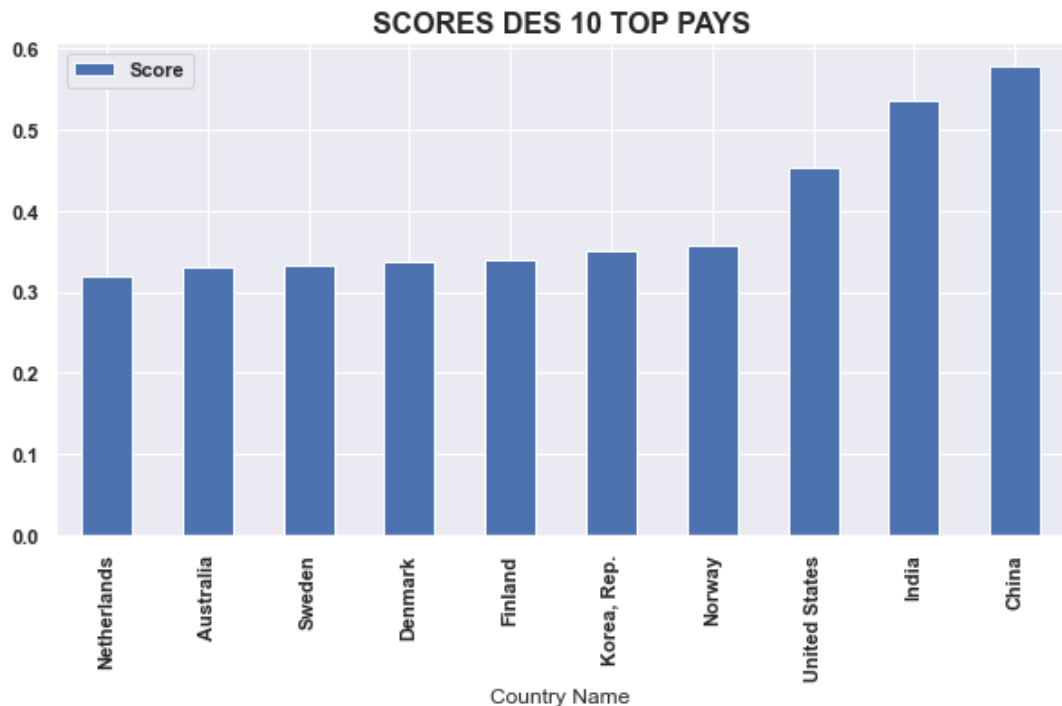
South Asia



Sub-Saharan Africa



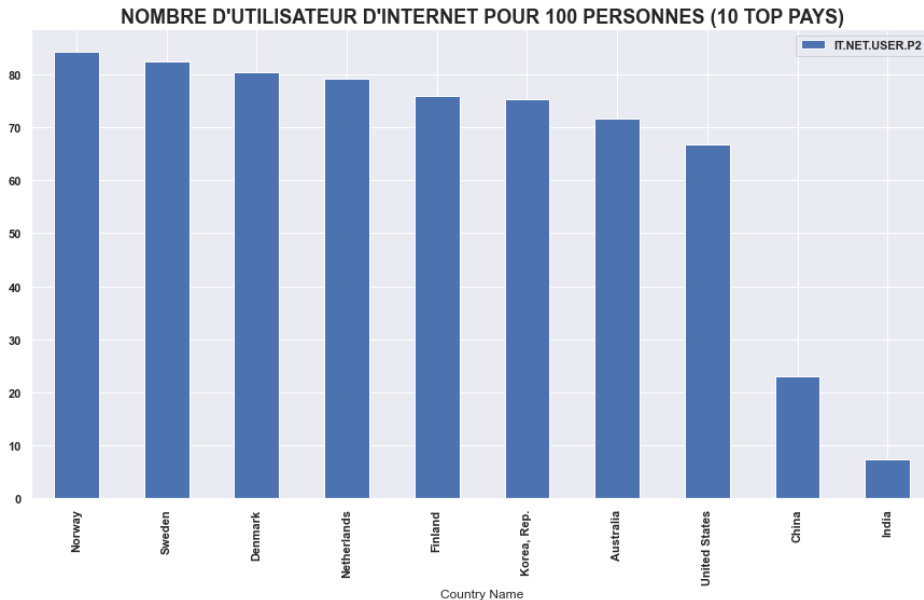
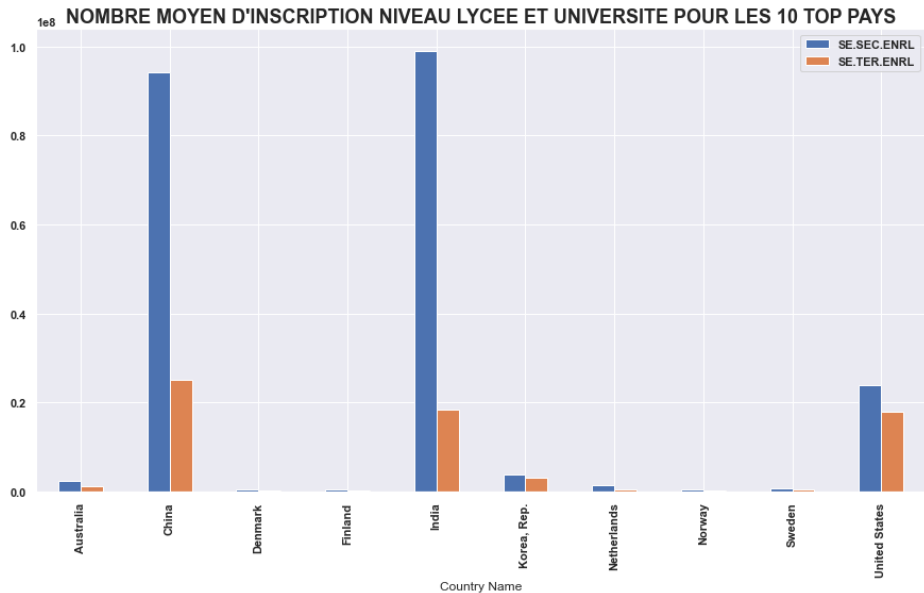
LE TOP 10 DES PAYS



OBSERVATIONS

- La Chine, l'Inde et les États-Unis sont en tête du classement
- Les données sont très bien remplies pour ces pays (ce qui crédibilise les scores obtenus)

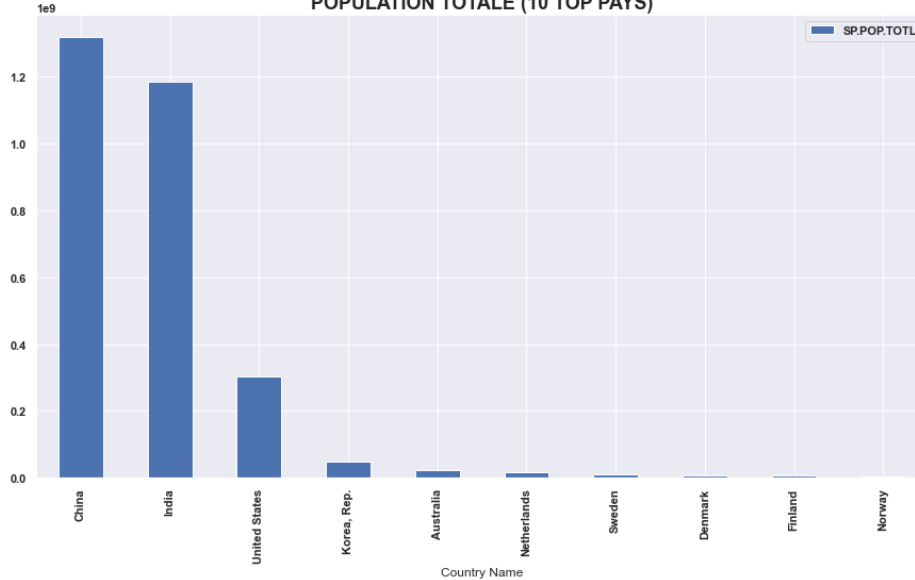
LYCÉE, UNIVERSITÉ, INTERNET



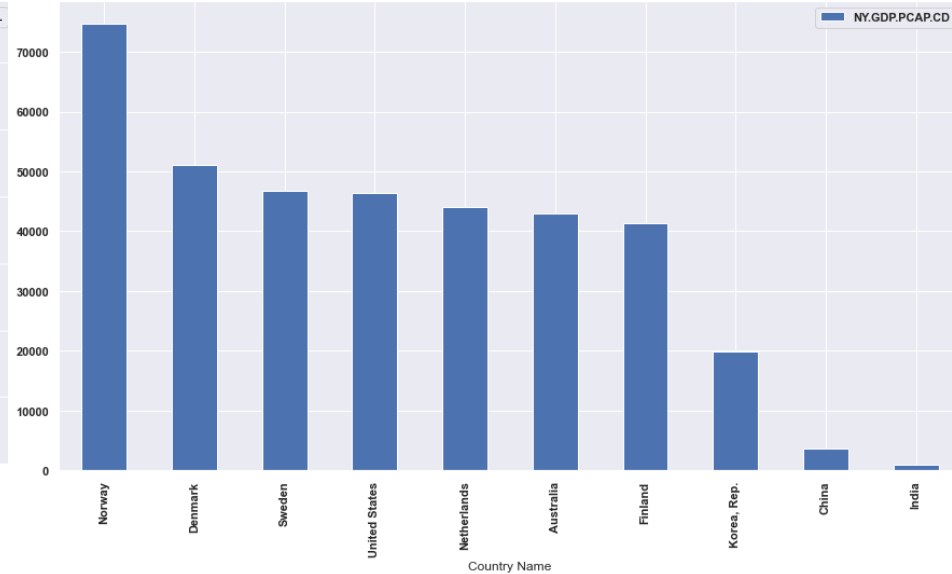
- La Chine et l'Inde ont un plus grand nombre de lycéens et d'étudiant
- Cependant ces 2 pays ont moins d'utilisateurs d'internet

POPULATION TOTALE, PIB/HABITANT

POPULATION TOTALE (10 TOP PAYS)



PIB/HABITANT (10 TOP PAYS)



- La Chine et l'Inde ont une population totale très élevée
- Cependant la Chine et l'Inde ont les plus bas PIB/habitant

Les scores élevés pour la Chine et l'Inde semblent être beaucoup plus liés aux indicateurs relatifs à la démographie

EVOLUTION DU POTENTIEL

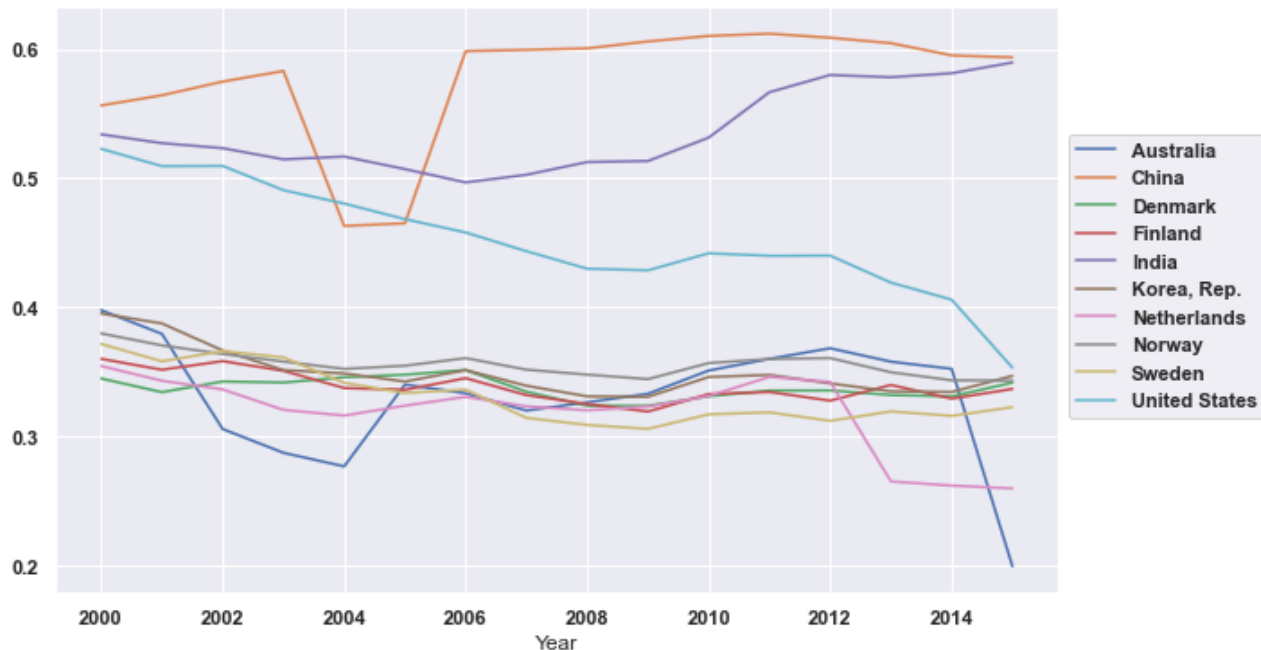
EVOLUTION DU SCORE DES 10 TOP PAYS

OBSERVATIONS

L'**Inde** semble être le pays ayant l'évolution la plus prometteuse

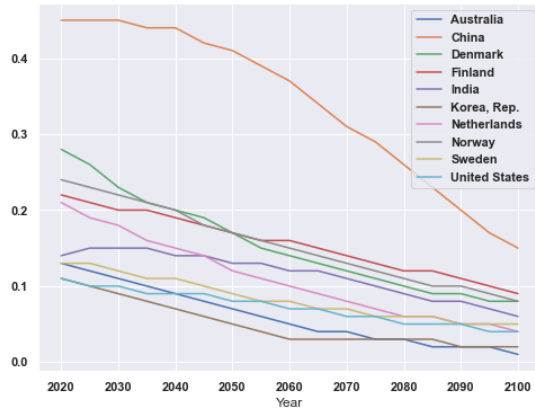
Tendance générale plutôt stable

A l'exception des Etats-Unis

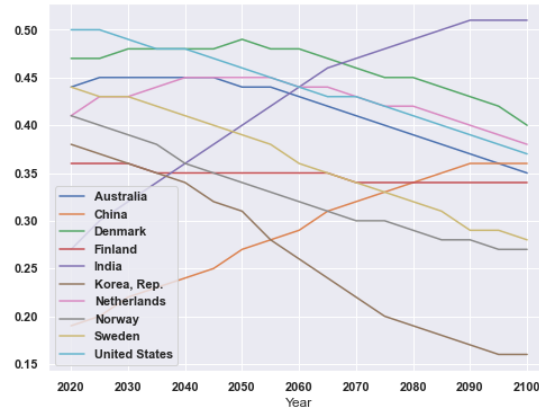


EVOLUTION INDICATEURS DE PROJECTION

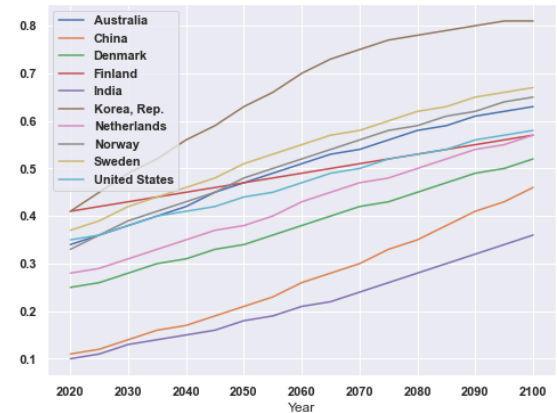
PROJECTION: EVOLUTION DU NIVEAU D'ETUDE (PRJ.ATT.15UP.2.MF)



PROJECTION: EVOLUTION DU NIVEAU D'ETUDE (PRJ.ATT.15UP.3.MF)



PROJECTION: EVOLUTION DU NIVEAU D'ETUDE (PRJ.ATT.15UP.4.MF)



Pour tous les 10 top pays, la tendance de l'évolution de la population de +15 ans ayant atteint le niveau "Post Secondary" sera croissante



03

CONCLUSION ET PERSPECTIVES

Suggestions de réponses aux questions
métier et perspectives d'amélioration

SUGGESTIONS

PAYS À FORT POTENTIEL

On pourrait sélectionner les pays avec les scores les plus élevés (Chine, Inde, Etats-Unis, Norvège, Corée, ...)

EVOLUTION DE CE POTENTIEL

La tendance du score au fil des années est plutôt stable pour la plupart des top pays excepté pour l'Inde qui semblent avoir une tendance d'évolution croissante

PAYS OÙ IL FAUDRAIT OPÉRER EN PRIORITÉ

Si l'on s'arrête à ce que nous avons proposé, l'entreprise devrait opérer en priorité dans le pays avec les scores les plus élevés. Par exemple: China, India, United States



PERSPECTIVES

- Sélection plus fine des indicateurs définitifs
- Attribution de poids différents aux indicateurs en fonction de leur niveau de pertinence d'un point de vue métier (par un expert du métier)

MERCI

Des questions ?

